# Allocations for Heterogenous Distributed Storage

Vasileios Ntranos
University of Southern California
Los Angeles, CA 90089, USA
ntranos@usc.edu

Giuseppe Caire
University of Southern California
Los Angeles, CA 90089, USA
caire@usc.edu

Alexandros G. Dimakis
University of Southern California
Los Angeles, CA 90089, USA
dimakis@usc.edu

*Abstract*—We study the problem of storing a data object in a set of data nodes that fail independently with given probabilities. Our problem is a natural generalization of a homogenous storage allocation problem where all the nodes had the same reliability and is naturally motivated for peer-to-peer and cloud storage systems with different types of nodes. Assuming optimal erasure coding (MDS), the goal is to find a storage allocation (i.e, how much to store in each node) to maximize the probability of successful recovery. This problem turns out to be a challenging combinatorial optimization problem. In this work we introduce an approximation framework based on large deviation inequalities and convex optimization. We propose two approximation algorithms and study the asymptotic performance of the resulting allocations. SUBMITTED TO ISIT 2012.

## I. INTRODUCTION

We are interested in heterogenous storage systems where storage nodes have different reliability parameters. This problem is relevant for heterogenous peer-to-peer storage networks and cloud storage systems that use multiple types of storage devices, *e.g.* solid state drives along with standard hard disks. We model this problem by considering $n$ storage nodes and a data collector that accesses a random subset $\mathbf{r}$ of them. The probability distribution of $\mathbf{r} \subseteq \{1, \ldots, n\}$ models random node failures and we assume that node $i$ fails independently with probability $1 - p_i$. The probability of a set $\mathbf{r}$ of nodes being accessed is therefore:

$$\mathbb{P}(\mathbf{r}) = \prod_{i \in \mathbf{r}} p_i \prod_{j \notin \mathbf{r}} (1 - p_j). \qquad (1)$$

Assume now that we have a single data file of unit size that we wish to code and store over these nodes to *maximize the probability of recovery* after a random set of nodes fail. The problem becomes trivial if we do not put a constraint on the maximum size $T$ of coded data and hence, we will work with a maximum storage budget of size $T < n$: If $x_i$ is the amount of coded data stored in node $i$, then $\sum_{i=1}^{n} x_i \leq T$. We further assume that our file is *optimally coded*, in the sense that successful recovery occurs whenever the total amount of data accessed by the data collector is at least the size of the original file. This is possible in practice when we use Maximum Distance Separable (MDS) codes [1]. The probability of successful recovery for an allocation $(x_1, \ldots, x_n)$ can be written as

$$P_s = \mathbb{P}\left[ \sum_{i \in \mathbf{r}} x_i \geq 1 \right] = \sum_{\mathbf{r} \subseteq \{1, \ldots, n\}} \mathbb{P}(\mathbf{r}) \, \mathbb{1}\left\{ \sum_{i \in \mathbf{r}} x_i \geq 1 \right\}$$

where $\mathbb{1}\{\cdot\}$ is the indicator function. $\mathbb{1}\{S\} = 1$ if the statement $S$ is true and zero otherwise.

A more concrete way to see this problem is by introducing a $Y_i \sim \text{Bernoulli}(p_i)$ random variable for each storage node: $Y_i = 1$ when node $i$ is accessed by the data collector and $Y_i = 0$ when node $i$ has failed. Define the random variable

$$Z = \sum_{i=1}^{n} x_i Y_i \qquad (2)$$

where $x_i$ is the amount of data stored in node $i$. Then, obviously, we have $P_s = \mathbb{P}[Z \geq 1]$.

Our goal is to find a storage allocation $(x_1, \ldots, x_n)$, that maximizes the probability of successful recovery, or equivalently, minimizes the probability of failure, $\mathbb{P}[Z < 1]$.

## II. OPTIMIZATION PROBLEM

Put in optimization form, we would like to find a solution to the following problem.

$$Q1: \quad \underset{x_i}{\text{minimize}} \quad \sum_{\mathbf{r} \subseteq \{1, \ldots, n\}} \mathbb{P}(\mathbf{r}) \, \mathbb{1}\left\{ \sum_{i \in \mathbf{r}} x_i < 1 \right\}$$

$$\text{subject to:} \quad \sum_{i=1}^{n} x_i \leq T$$

$$x_i \geq 0, \ i = 1, \ldots, n.$$

Authors in [1] consider a special case of problem $Q1$ in which $p_i = p$, $\forall i$. Even in this symmetric case the problem appears to be very difficult to solve due to its non-convex and combinatorial nature. In fact, even for a given allocation $\{x_i\}$ and parameter $p$, computing the objective function is computationally intractable ($\#P\text{-}hard$, See [1]).

A very interesting observation about this problem follows directly from Markov's Inequality: $\mathbb{P}[Z \geq 1] \leq E[Z] = pT$. If $pT < 1$, then the probability of successful recovery is bounded away from 1. This has motivated the definition of a region of parameters for which high probability of recovery is possible: $R_{HP} = \{(p, T) : pT \geq 1\}$. The budget $T$ should be more than $1/p$ if we want to aim for high reliability and the authors in [1] showed that in the above region of parameters, *maximally spreading* the budget to all nodes (i.e, $x_i = T/n$, $\forall i$) is an asymptotically optimal allocation as $n \to \infty$.

In the general case, when the node access probabilities, $p_i$, are not equal, one could follow similar steps to characterize

a region of high probability of recovery. Markov's Inequality yields:

$$\mathbb{P}[Z \geq 1] \leq E[Z] = \sum_{i=1}^{n} x_i p_i = \mathbf{p}^T \mathbf{x}$$

where $\mathbf{p} = [p_1, p_2, \ldots, p_n]^T$ and $\mathbf{x} = [x_1, x_2, \ldots, x_n]^T$. If we don't want $\mathbb{P}[Z \geq 1]$ to be bounded away from 1 we have to require now that $\mathbf{p}^T \mathbf{x} \geq 1$. We see that in this case, high reliability is not a matter of sufficient budget, as it depends on the allocation $\mathbf{x}$ itself.

Let $S(\mathbf{p}, T) = \left\{ \mathbf{x} \in \mathbb{R}_+^n : \mathbf{p}^T \mathbf{x} \geq 1, \mathbf{1}^T \mathbf{x} \leq T \right\}$ be the set of all allocations $\mathbf{x}$ with a given budget constraint $T$ that satisfy $\mathbf{p}^T \mathbf{x} \geq 1$ for a given $\mathbf{p}$. We call these allocations *reliable* for a system with parameters $\mathbf{p}, T$, in the sense that the resulting probability of successful recovery is not bounded away from 1. Then the region of high probability of recovery can be defined as the region of parameters $\mathbf{p}, T$, such that the set $S(\mathbf{p}, T)$ is non-empty.

$$\mathcal{R}_{HP} = \left\{ (\mathbf{p}, T) \in \mathbb{R}_+^{n+1} : S(\mathbf{p}, T) \neq \emptyset \right\}$$

This generalizes the region described in [1]. If all $p_i$'s are equal then the set $S(\mathbf{p}, T)$ is non-empty when $\mathbf{p}^T \mathbf{x} = pT \geq 1$. In the general case, the minimum budget such that $S(\mathbf{p}, T)$ is non-empty is $T = 1/p_{max}$, with $p_{max} = \max\{p_i\}$, and $S(\mathbf{p}, 1/p_{max})$ contains only one allocation $\mathbf{x}_{p_{max}^{-1}} : x_j = \frac{1}{p_{max}}$, $j = arg \max_i \{p_i\}$, $x_i = 0$, $\forall i \neq j$.

Even though $\mathcal{R}_{HP}$ provides a lower bound on the minimum budget $T$ required to allocate for high reliability, it doesn't provide any insights on how to design allocations that achieve high probability of recovery in a distributed storage system. This motivates us to move one step further and define a region of $\epsilon$-optimal allocations in the next section.

## III. THE REGION OF $\epsilon$-OPTIMAL ALLOCATIONS

We say that an allocation $(x_1, x_2, \ldots, x_n)$ is $\epsilon$-optimal if the corresponding probability of successful recovery, $\mathbb{P}[Z \geq 1]$, is greater than $1 - \epsilon$.

Let $\mathcal{E}_n(\mathbf{p}, T, \epsilon) = \{ \mathbf{x} \in \mathbb{R}_+^n : \mathbb{P}[Z < 1] \leq \epsilon, \mathbf{1}^T \mathbf{x} \leq T \}$ be the set of all $\epsilon$-optimal allocations. Note that if we could *efficiently* characterize this set for all problem parameters, we would be able to solve problem $Q1$ exactly: Find the smallest $\epsilon$ such that $\mathcal{E}_n(\mathbf{p}, T, \epsilon)$ is non-empty.

In this section we will derive a sufficient condition for an allocation to be $\epsilon$-optimal and provide an efficient characterization for a region $\mathcal{H}_n \subseteq \mathcal{E}_n(\mathbf{p}, T, \epsilon)$. We begin with a very useful lemma.

**Lemma 1.** *(Hoeffding's Inequality [2], [3])*
*Consider the random variable $W = \sum_{i=1}^{n} V_i$, where $V_i$ are independent almost surely bounded random variables with $\mathbb{P}(V_i \in [a_i, b_i]) = 1$. Then,*

$$\mathbb{P}\left[ W \leq E[W] - n\delta \right] \leq exp\left\{ -\frac{2n^2\delta^2}{\sum_{i=1}^{n}(b_i - a_i)^2} \right\}$$

*for any $\delta > 0$.*

We can use Lemma 1 to upper bound the probability of failure, $\mathbb{P}[Z < 1] \leq \mathbb{P}[Z \leq 1]$, for an arbitrary allocation, since $Z = \sum_{i=1}^{n} x_i Y_i$ can be seen as the sum of $n$ independent almost surely bounded random variables $V_i = x_i Y_i$, with $\mathbb{P}(V_i \in [0, x_i]) = 1$. Let $\delta = \left( \sum_{i=1}^{n} x_i p_i - 1 \right)/n$ and require $\delta > 0 \Leftrightarrow \sum_{i=1}^{n} x_i p_i > 1$. Lemma 1 yields:

$$\mathbb{P}[Z < 1] \leq exp\left\{ -\frac{2\left( \sum_{i=1}^{n} x_i p_i - 1 \right)^2}{\sum_{i=1}^{n} x_i^2} \right\}, \quad \sum_{i=1}^{n} x_i p_i > 1. \tag{3}$$

Notice that the constraint $\sum_{i=1}^{n} x_i p_i > 1$ requires the allocation $(x_1, x_2, \ldots, x_n)$ to be reliable and $S(\mathbf{p}, T) \neq \emptyset$.

In view of the above, a sufficient condition for a strictly reliable allocation to be $\epsilon$-optimal is the following.

$$exp\left\{ -\frac{2\left( \sum_{i=1}^{n} x_i p_i - 1 \right)^2}{\sum_{i=1}^{n} x_i^2} \right\} \leq \epsilon \quad \Longleftrightarrow$$
$$||\mathbf{x}||_2 \sqrt{\frac{\ln 1/\epsilon}{2}} \leq \mathbf{p}^T \mathbf{x} - 1, \quad \mathbf{p}^T \mathbf{x} > 1 \tag{4}$$

We say that all allocations satisfying the above equation are *Hoeffding $\epsilon$-optimal*, due to the use of Hoeffding's Inequality in Lemma 1.

**Definition 1.** *"The Region of Hoeffding $\epsilon$-optimal allocations"*

$$\mathcal{H}_n(\mathbf{p}, T, \epsilon) = \left\{ \mathbf{x} \in \mathbb{R}_+^n : \mathbf{p}^T \mathbf{x} > 1, \mathbf{1}^T \mathbf{x} \leq T, \right.$$
$$\left. ||\mathbf{x}||_2 \sqrt{\frac{\ln 1/\epsilon}{2}} \leq \mathbf{p}^T \mathbf{x} - 1 \right\} \tag{5}$$

The above region is strictly smaller $\mathcal{E}_n(\mathbf{p}, T, \epsilon)$ for any finite $n$, because the bound in (3) is not generally tight. However, $\mathcal{H}_n(\mathbf{p}, T, \epsilon)$ is a convex set: Equation (4) can be seen as a second order cone constraint on the allocation $\mathbf{x} \in \mathbb{R}_+^n$.

**Theorem 1.** *The region of Hoeffding $\epsilon$-optimal allocations $\mathcal{H}_n(\mathbf{p}, T, \epsilon)$ is convex in $\mathbf{x}$.*

This interesting result allows us to formulate and efficiently solve optimization problems over $\mathcal{H}_n(\mathbf{p}, T, \epsilon)$. Finding the smallest $\epsilon^*$ such that $\mathcal{H}_n(\mathbf{p}, T, \epsilon)$ is non-empty will produce an $\epsilon^*$-optimal solution to problem $Q1$.

### A. Hoeffding Approximation of $Q1$

If we fix $\mathbf{p}, T, n$ as the problem parameters, then the following optimization problem can be solved efficiently, to any desired accuracy $1/\alpha$, by solving a sequence of $\mathcal{O}(\log \alpha)$ convex feasibility problems (bisection on $\epsilon$).

$$H1: \quad \min_{\mathbf{x}, \epsilon} \quad \epsilon$$
$$\text{s.t.:} \quad \mathbf{x} \in \mathcal{H}_n(\mathbf{p}, T, \epsilon)$$

We will see next that if $T$ is sufficiently large, $\epsilon^*$ goes to zero exponentially fast as $n$ grows, and hence the solution to the aforementioned problem is asymptotically optimal.

## B. Maximal Spreading Allocations and the Asymptotic Optimality of H1

First, we will focus on maximal spreading allocations, $\mathbf{x}_T^n \triangleq \{\mathbf{x} \in \mathbb{R}^n : x_i = T/n\}$, and derive their asymptotic optimality for $Q1$, in the sense that $\mathbb{P}[Z < 1] \to 0$, as $n \to \infty$. Let $\bar{p} = \frac{1}{n}\sum_{i=1}^n p_i$ be the average access probability across all nodes. We have the following lemma.

**Lemma 2.** *If $T > 1/\bar{p}$, for any $\epsilon > 0$, $\exists n_\epsilon$: $\mathbf{x}_T^n \in \mathcal{H}_n(\mathbf{p}, T, \epsilon)$, for all $n \geq n_\epsilon$.*

*Proof:* This follows directly from the definition of $\mathcal{H}_n(\mathbf{p}, T, \epsilon)$: $n_\epsilon = \frac{\ln 1/\epsilon}{2(\bar{p}-1/T)^2}$ . ∎

The above lemma establishes the asymptotic optimality of maximal spreading allocations through the following corollary.

**Corollary 1.** *The probability of failed recovery, $P_e \triangleq \mathbb{P}[Z < 1]$, for a maximal spreading allocation is $P_e \leq e^{-2n(\bar{p}-1/T)^2}$. When $T > 1/\bar{p}$, $P_e \to 0$, as $n \to \infty$.*

The fact that $\mathcal{H}_n(\mathbf{p}, T, \epsilon)$ contains maximal spreading allocations for $T > 1/\bar{p}$, provides a sufficient condition on the asymptotic optimality of $H1$.

**Theorem 2.** *Let $\epsilon^*$ be the optimal value of $H1$. If $T > 1/\bar{p}$, then $\epsilon^* = \mathcal{O}(exp(-n))$.*

*Proof:* Let $T > 1/\bar{p}$ and consider the maximal spreading allocation $\mathbf{x}_T^n$. Then, $\epsilon^* \leq \epsilon_s$, where $\epsilon_s$ is the minimum $\epsilon$ such that $\mathbf{x}_T^n \in \mathcal{H}_n(\mathbf{p}, T, \epsilon)$. That is $\epsilon_s = e^{-2n(\bar{p}-1/T)^2}$, and since $T > 1/\bar{p}$, $\epsilon^* \leq \epsilon_s = \mathcal{O}(exp(-n))$. ∎

## IV. CHERNOFF RELAXATION

In this section we take a different approach to obtain a tractable convex relaxation for $Q1$ by minimizing an appropriate Chernoff upper bound.

### A. Upper Bounding the Objective Function

**Lemma 3.** *(Upper Bound) Let $Z = \sum_{i=1}^n x_i Y_i$, $x_i \geq 0$, $Y_i \sim bernoulli(p_i)$ and $t \geq 0$. The probability of failed recovery, $\mathbb{P}[Z < 1]$, is upper bounded by*

$$\mathbb{P}[Z < 1] \leq g_t(\mathbf{x}) = \sum_{\mathbf{r} \subseteq \{1,\ldots,n\}} \mathbb{P}(\mathbf{r}) \exp\left\{-t\left(\sum_{i \in \mathbf{r}} x_i - 1\right)\right\}$$

*Proof:* For any $t \geq 0$ we have:

$$\begin{aligned}
\mathbb{P}[Z < 1] &\leq \mathbb{P}[Z \leq 1] = \mathbb{P}\left[e^{-tZ} \geq e^{-t}\right] \\
&\leq e^t \mathbb{E}\left[e^{-tZ}\right] = e^t \mathbb{E}\left[\prod_{i=1}^n e^{-tx_i Y_i}\right] \\
&= e^t \prod_{i=1}^n \mathbb{E}\left[e^{-tx_i Y_i}\right] \\
&= e^t \prod_{i=1}^n \left(1 - p_i + p_i e^{-tx_i}\right) \quad (6)
\end{aligned}$$

$$\begin{aligned}
&= e^t \sum_{\mathbf{r} \subseteq \{1,\ldots,n\}} \mathbb{P}(\mathbf{r}) \exp\left\{-t\left(\sum_{i \in \mathbf{r}} x_i\right)\right\} \\
&= \sum_{\mathbf{r} \subseteq \{1,\ldots,n\}} \mathbb{P}(\mathbf{r}) \exp\left\{-t\left(\sum_{i \in \mathbf{r}} x_i - 1\right)\right\} \quad (7) \\
&\triangleq g_t(\mathbf{x})
\end{aligned}$$

∎

Note that $g_t(\mathbf{x})$ is a weighted sum of convex functions with linear arguments, and hence convex in $\mathbf{x}$. Equation (7) makes the convex relaxation of the objective function apparent:

$$\mathbb{1}\left\{x < \alpha\right\} \leq e^{-t(x-\alpha)}, \text{ for any } t \geq 0.$$

### B. The Relaxed Optimization Problem

Before we move forward and state the relaxed optimization problem, we take a closer look at the constraint set $S = \{\mathbf{x} \in \mathbb{R}_+^n : \mathbf{1}^T \mathbf{x} \leq T\}$ of the original problem $Q1$. From a practical perspective, it should be wasteful to allocate more than one unit of data (filesize) on a single node. If the node survives, then the data collector can always recover the file using only one unit of data and hence any additional storage does not help. Also, an allocation using less than the available budget cannot have larger probability of successful recovery.

In the following lemma, we show that it is sufficient to consider allocations with $x_i \in [0, 1]$ and $\sum_{i=1}^n x_i = T$.

**Lemma 4.** *For any $\mathbf{x} \in S$, $\exists \mathbf{x}' \in S' = \{\mathbf{x} \in \mathbb{R}_+^n : \mathbf{1}^T \mathbf{x} = T, x_i \leq 1, i = 1, \ldots, n\}$ such that $\mathbb{P}\left[\sum_{i=1}^n x_i' Y_i < 1\right] \leq \mathbb{P}\left[\sum_{i=1}^n x_i Y_i < 1\right]$.*

*Proof:* See the long version of this paper [4]. ∎

The relaxed optimization problem can be formulated as follows.

$$\begin{aligned}
R1: \quad &\underset{x_i}{\text{minimize}} \quad g_t(\mathbf{x}) \\
&\text{subject to:} \quad \sum_{i=1}^n x_i = T \\
&\qquad\qquad\quad x_i \in [0, 1], \; i = 1, \ldots, n.
\end{aligned}$$

Note that, in general, one would like to minimize $\inf_{t \geq 0}\{g_t(\mathbf{x})\}$ instead of $g_t(\mathbf{x})$ for some $t \geq 0$. However, for now, we will let $t$ be a free parameter and carry on with the optimization.

The important drawback of the above formulation hides in the objective function: Although convex, $g_t(\mathbf{x})$ has an *exponentially long* description in the number of storage nodes: The sum is still over all subsets $\mathbf{r} \subseteq \{1, \ldots, n\}$. This can be circumvented if we consider minimizing $\log g_t(\mathbf{x})$ instead of $g_t(\mathbf{x})$ over the same set.

**Lemma 5.** *$\log g_t(\mathbf{x})$ is convex in $\mathbf{x}$.*

*Proof:* See the long version of this paper [4].

∎

**Lemma 6.** *For any $t \geq 0$*

$$\arg\min_{\mathbf{x} \in S} g_t(\mathbf{x}) = \arg\min_{\mathbf{x} \in S} \sum_{i=1}^{n} \log\left(1 + \frac{p_i}{1 - p_i} e^{-tx_i}\right),$$

*where $S = \{\mathbf{x} \in \mathbb{R}_+^n : \mathbf{1}^T \mathbf{x} \leq T, \mathbf{x} \preceq \mathbf{1}\}$.*

*Proof:* Let $\mathbf{x}^* = \arg\min_{\mathbf{x} \in S} g_t(\mathbf{x})$. Then $g_t(\mathbf{x}^*) \leq g_t(\mathbf{x})$, $\forall \mathbf{x} \in S$. Taking the logarithm on both sides preserves the inequality since $\log(\cdot)$ is strictly increasing. Hence, $\log g_t(\mathbf{x}^*) \leq \log g_t(\mathbf{x})$, $\forall \mathbf{x} \in S$ and subtracting $t + \sum_{i=1}^{n} \log(1 - p_i)$ from both sides yields the desired result and completes the proof. ∎

In view of Lemmas 5 and 6, we can solve $R1$ through the following *equivalent* optimization problem.

$$R2: \quad \underset{x_i}{\text{minimize}} \quad t + \sum_{i=1}^{n} \log\left(1 + \frac{p_i}{1 - p_i} e^{-tx_i}\right)$$

$$\text{subject to:} \quad \sum_{i=1}^{n} x_i = T$$

$$x_i \in [0, 1], \ i = 1, \ldots, n.$$

$R2$ is a convex *separable* optimization problem with polynomial size description and in terms of complexity, it is "not much harder" than linear programming [5]. One can solve such problems numerically in a very efficient way using standard, "off-the-shelf" algorithms and optimization packages such as CVX [6], [7].

### C. Insights from Optimality Conditions for R2

Here, we move one step further and take the KKT conditions for $R2$ in order to take a closer look at the structure of the optimal solutions. Let $r_i \triangleq \frac{p_i}{1 - p_i}$.

The Lagrangian for $R2$ is:

$$L(\mathbf{x}, \mathbf{u}, \mathbf{v}, \lambda) = \sum_{i=1}^{n} \log\left(1 + r_i e^{-tx_i}\right) + \lambda\left(\sum_{i=1}^{n} x_i - T\right)$$
$$- \sum_{i=1}^{n} u_i x_i + \sum_{i=1}^{n} v_i(x_i - 1)$$

where $\lambda \in \mathbb{R}$, $\mathbf{u}, \mathbf{v} \in \mathbb{R}_+^n$ are the corresponding Lagrange multipliers. The gradient is given by $\nabla_{x_i} L(\mathbf{x}, \mathbf{u}, \mathbf{v}, \lambda) = -\frac{r_i t}{r_i + e^{tx_i}} + \lambda - u_i + v_i$, and the KKT necessary and sufficient conditions for optimality yield:

$$-\frac{r_i t}{r_i + e^{tx_i^*}} + \lambda - u_i + v_i = 0, \ \forall i \tag{8}$$

$$\sum_{i=1}^{n} x_i^* = T \tag{9}$$

$$0 \leq x_i^* \leq 1, \ \forall i \tag{10}$$

$$\lambda \in \mathbb{R}, \ v_i, u_i \geq 0, \ \forall i \tag{11}$$

$$v_i(x_i - 1) = 0, \ u_i x_i = 0, \ \forall i \tag{12}$$

Using the results from [8], the optimal solution to $R2$ is given by

$$x_i^* = \begin{cases} 0 & \text{if } \frac{r_i t}{1 + r_i} \leq \lambda^* \\ 1 & \text{if } \lambda^* \leq \frac{r_i t}{e^t + r_i} \\ \frac{1}{t} \log\left(\frac{r_i t}{\lambda^*} - r_i\right) & \text{if } \frac{r_i t}{e^t + r_i} < \lambda^* < \frac{r_i t}{1 + r_i} \end{cases} \tag{13}$$

where $\lambda^*$ is chosen such that Eq.(9) is satisfied, i.e,

$$\sum_{i=1}^{n} \frac{1}{t} \log\left(\frac{r_i t}{\lambda^*} - r_i\right) \mathbb{1}\left\{\lambda^* \in \left(\frac{r_i t}{e^t + r_i}, \frac{r_i t}{1 + r_i}\right)\right\}$$
$$+ \sum_{i=1}^{n} \mathbb{1}\left\{\lambda^* \leq \frac{r_i t}{e^t + r_i}\right\} = T \tag{14}$$

Numerically, $\lambda^*$ can be computed via an iterative $\mathcal{O}(n^2)$ algorithm described in [8], and hence this approach gives an even more efficient way to solve $R2$.

However, the most important aspect of the above result is that we can use equations (13), (14) to obtain closed form solutions for a certain region of problem parameters and analyze the performance of the resulting allocations.

### D. The choice of parameter $t \geq 0$

It is clear that the optimal solution to $R2$ depends on our choice of $t \geq 0$. For example, $\frac{r_i t}{e^t + r_i} \to 0$, $\frac{r_i t}{1 + r_i} \to \infty$, as $t \to \infty$ and $x_i^* = \lim_{t \to \infty} t^{-1} \log\left(r_i t/\lambda^* - r_i\right)$, $\forall i$. Equation (14) yields $x_i^* = \frac{T}{n}$, $\forall i$ and hence the maximal spreading allocation becomes optimal for $R2$ as $t \to \infty$. Even though this motivates the choice of maximal spreading allocations as approximate "one-shot" solutions for the original problem $Q1$, explicitly tuning the parameter $t$ can provide significantly better approximations.

In order to obtain the tightest bound from Lemma 3, we have to jointly minimize the objective in $R2$ with respect to $t \geq 0$ and $\mathbf{x}$. Towards this end, one can *iteratively optimize* $R2$ by fixing the value of one variable ($t$ or $\mathbf{x}$) at each step and minimizing over the other. After each iteration the objective function decreases and hence the above procedure converges to a (possibly local) minimum. The above algorithm iteratively tunes the Chernoff bound introduced in this section and produces a minimizing allocation that can serve as an approximate solution to the original problem $Q1$.

For analytic purposes though, we can choose a value for $t$ as follows. Recall from Lemma 3 that $\mathbb{P}[Z < 1] \leq g_t(\mathbf{x})$ for any $t \geq 0$. After taking logarithms, we would like to find a value for $t \geq 0$ that minimizes $b(t) \triangleq t + \sum_{i=1}^{n} \log(1 + r_i e^{-tx_i})$. Notice that $b(t)$ is a convex function of $t$, with $b(t) > 0$, $\forall t \geq 0$, $b(0) = \sum_{i=1}^{n} \log(1 + r_i)$ and $\lim_{t \to \infty} b(t) = \infty$. The slope of $b(t)$ at zero is $b'(0) = 1 - \sum_{i=1}^{n} \frac{r_i x_i}{1 + r_i} = 1 - \sum_{i=1}^{n} p_i x_i$, which is negative if the allocation is reliable.

When $t$ is large, $\log(1 + r_i e^{-tx_i}) \approx 0$, whereas for small values of $t$, $\log(1 + r_i e^{-tx_i}) \approx -tx_i + \log r_i$ and hence $b(t) \approx t + \sum_{i=1}^{n} \max\{-tx_i + \log r_i, 0\} \geq t + \max\{-\sum_{i=1}^{n} tx_i + \log r_i, 0\}$. One way to choose $t$ that does not depend on $x_i$ is to make $-\sum_{i=1}^{n} tx_i + \log r_i = 0 \Rightarrow t = \frac{1}{T} \sum_{i=1}^{n} \log r_i$.

## E. A closed-form allocation: $\hat{\mathbf{x}}_T^n$

In view of the above results we provide here a closed form allocation (each $x_i$ is given as a function of $\mathbf{p}$ and T) that can be used to study the asymptotic performance of $R2$ and serve as a better "one-shot" approximate solution to Q1.

Let $\mathcal{E}(\cdot)$ be a shorthand notation for the sample average such that $\mathcal{E}f(x) = \frac{1}{n}\sum_{i=1}^{n} f(x_i)$, in order to simplify the expressions. For the above choice of $t = \frac{1}{T}\sum_{i=1}^{n}\log r_i = n\mathcal{E}\log r/T$, equation (13) becomes:

$$x_i^* = \begin{cases} 0 & \text{if } \frac{r_i}{1+r_i}\frac{n\mathcal{E}\log r}{T} \leq \lambda^* \\ 1 & \text{if } \lambda^* \leq \frac{r_i n\mathcal{E}\log r/T}{e^{n\mathcal{E}\log r/T}+r_i} \\ \frac{T}{n\mathcal{E}\log r}\log\left(\frac{nr_i\mathcal{E}\log r}{T\lambda^*} - r_i\right) & \text{otherwise} \end{cases}$$

(15)

**Lemma 7.** If $p_i > \frac{1}{2}$, $\forall i$ and $T < \frac{n\mathcal{E}\log r}{\log r_{max}}$, $r_{max} = max\{r_i\}$, then $x_i^* = \frac{T}{n\mathcal{E}\log r}\log r_i$, $\forall i$.

*Proof:* Assume that $\lambda^* \in \left(\frac{r_i n\mathcal{E}\log r/T}{e^{n\mathcal{E}\log r/T}+r_i}, \frac{r_i}{1+r_i}\frac{n\mathcal{E}\log r}{T}\right)$. Then from Eq.(14), $\lambda^* = \frac{n\mathcal{E}\log r}{2T}$ and $x_i^* = \frac{T}{n\mathcal{E}\log r}\log r_i$. $\lambda^*$ is indeed in the required interval if $\frac{n\mathcal{E}\log r}{2T} < \frac{r_i}{1+r_i}\frac{n\mathcal{E}\log r}{T}$, $\forall i$ $\Rightarrow r_i > 1$, $\forall i \Rightarrow p_i > 1/2$, $\forall i$ and $\frac{n\mathcal{E}\log r}{2T} < \frac{r_i n\mathcal{E}\log r/T}{e^{n\mathcal{E}\log r/T}+r_i}$, $\forall i$ $\Rightarrow r_i < e^{n\mathcal{E}\log r}/T$, $\forall i \Rightarrow T < \frac{n\mathcal{E}\log r}{\log r_{max}}$. ∎

Clearly, when all $p_i > 1/2$, $\hat{\mathbf{x}}_T^n$: $x_i = \frac{T}{n\mathcal{E}\log r}\log r_i$, $\forall i$, is a feasible suboptimal allocation for $Q1$. It is also suboptimal for $R2$ in general, since solving $R2$ via the proposed algorithms can only achieve a smaller probability of failed recovery. We have $P_e\{Q1\} \leq P_e\{R2\} \leq \mathbb{P}\left\{\sum_{i=1}^{n}\frac{T}{n\mathcal{E}\log r}\log r_i Y_i < 1\right\}$.

In the following lemma we give an upper bound on the probability of failed recovery for $\hat{\mathbf{x}}_T^n$ and establish its asymptotic optimality.

**Lemma 8.** If $p_i > \frac{1}{2}$, $\forall i$ and $T > \frac{\mathcal{E}\log r}{\mathcal{E}p\log r}$, the allocation $\hat{\mathbf{x}}_T^n$ : $x_i = \frac{T}{n\mathcal{E}\log r}\log r_i$, $\forall i$, is strictly reliable, and the probability of failed recovery, $P_e = \mathbb{P}[Z < 1]$, is upper bounded by

$$P_e \leq exp\left\{-2n\frac{\left(\mathcal{E}p\log r - \frac{\mathcal{E}\log r}{T}\right)^2}{\mathcal{E}\log^2 r}\right\}$$

and hence, when $T > \frac{\mathcal{E}\log r}{\mathcal{E}p\log r}$, $P_e \to 0$, as $n \to \infty$.

*Proof:*
The proof follows directly from Lemma 1 and Equation (3). ∎

Notice that $\hat{\mathbf{x}}_T^n$ is reliable for values of $T$ for which a maximal spreading allocation $\mathbf{x}_T^n$ is not, since $\frac{1}{\bar{p}} \geq \frac{\mathcal{E}\log r}{\mathcal{E}p\log r}$, and hence its probability of failed recovery $P_e$ goes to zero exponentially fast for smaller values of $T$.

## V. NUMERICAL EXPERIMENTS

In this section we evaluate the performance of the proposed approximate distributed storage allocations in terms of their probability of failed recovery and plot the corresponding
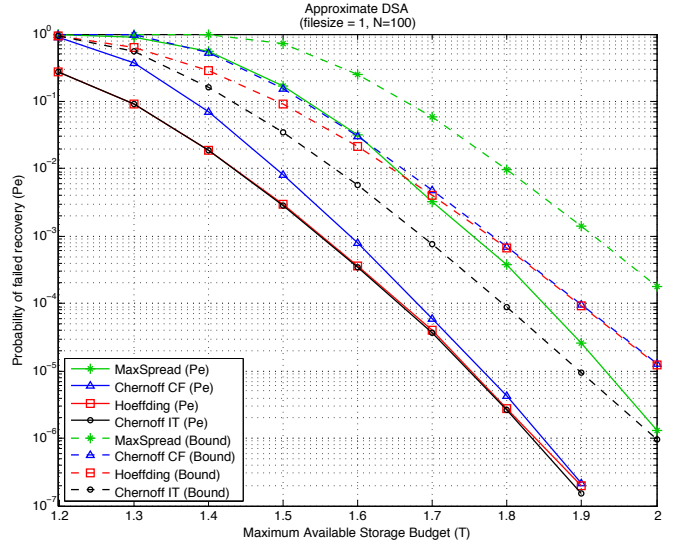


Fig. 1. Performance of the proposed approximate distributed storage allocations and their corresponding upper bounds for a system with $n = 100$ nodes and $p_i \sim \mathcal{U}(0.5, 1)$.

bounds. In our simulations we consider an ensemble of distributed storage systems with $n = 100$ nodes, in which the corresponding access probabilities, $p_i \sim \mathcal{U}(0.5, 1)$, are drawn uniformly at random from the interval $(0.5, 1)$.

We consider the following allocations. 1) *Maximal spreading*: $x_i = \frac{T}{n}$, $\forall i$. 2) *Chernoff closed-form*: $x_i = (T/n\mathcal{E}\log r)\log r_i$, $\forall i$. 3) *Hoeffding $\epsilon$-optimal*: obtained by solving $H1$. 4) *Chernoff iterative*: obtained by solving $R2$ and iteratively tuning the parameter $t$.

Fig.1 shows, in solid lines, the ensemble average probability of failed recovery of each allocation, $\mathbb{P}\left[\sum_{i=1}^{n} x_i Y_i < 1\right]$, versus the maximum available budget $T$. In dashed lines, Fig.1 plots the corresponding bounds on $P_e$ obtained from Corollary 1, Lemma 8 and the objective functions of $H1$, $R1$.

## REFERENCES

[1] D. Leong, A. Dimakis, and T. Ho. Distributed storage allocations. *CoRR*, abs/1011.5287, 2010.
[2] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Stat. Association*, 58(301):13–30, March 1963.
[3] M. Mitzenmacher and E. Upfal. *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, New York, NY, USA, 2005.
[4] V. Ntranos, G. Caire, and A. Dimakis. Allocations for heterogenous distributed storage (*long version*). http://www-scf.usc.edu/~ntranos/docs/HDS-long.pdf, January 2012.
[5] D. S. Hochbaum and J. George Shanthikumar. Convex separable optimization is not much harder than linear optimization. *J. ACM*, 37:843–862, October 1990.
[6] M. Grant and S. Boyd. CVX: Matlab software for disciplined convex programming, version 1.21. http://cvxr.com/cvx, April 2011.
[7] M. Grant and S. Boyd. Graph implementations for nonsmooth convex programs. In V. Blondel, S. Boyd, and H. Kimura, editors, *Recent Advances in Learning and Control*, Lecture Notes in Control and Information Sciences, pages 95–110. Springer-Verlag Limited, 2008.
[8] S. M. Stefanov. Convex separable minimization subject to bounded variables. *Comp. Optimization and Applications*, 18, 2001.